

Computational & Information Systems Laboratory

Performance Analysis of MPI over InfiniBand on Yellowstone

Zhengyang Liu
Mentor: Dr. John Dennis

Collaborators:
Prof. Malathi Veeraraghavan (University of Virginia)
Prof. Robert D. Russell (University of New Hampshire)
Fabrice Mizero (SIParCS)
Patrick MacArthur (University of New Hampshire)

Oct 25, 2013

NCAR

UNIVERSITY of VIRGINIA

Computational & Information Systems Laboratory

Big Picture

- **Understanding the causes of poor performance of CESM on Yellowstone: a 5-step approach**
 - Experimental execution and data collection
 - CESM trace analysis
 - IBMgtSim: routing study
 - Network simulation
 - Integrated simulation

NCAR

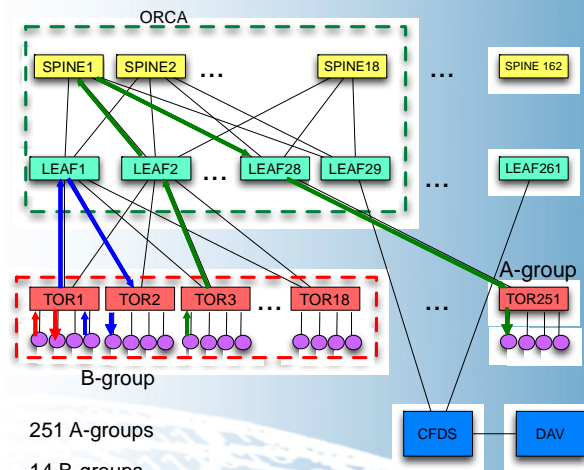
UNIVERSITY of VIRGINIA

1

Big Picture

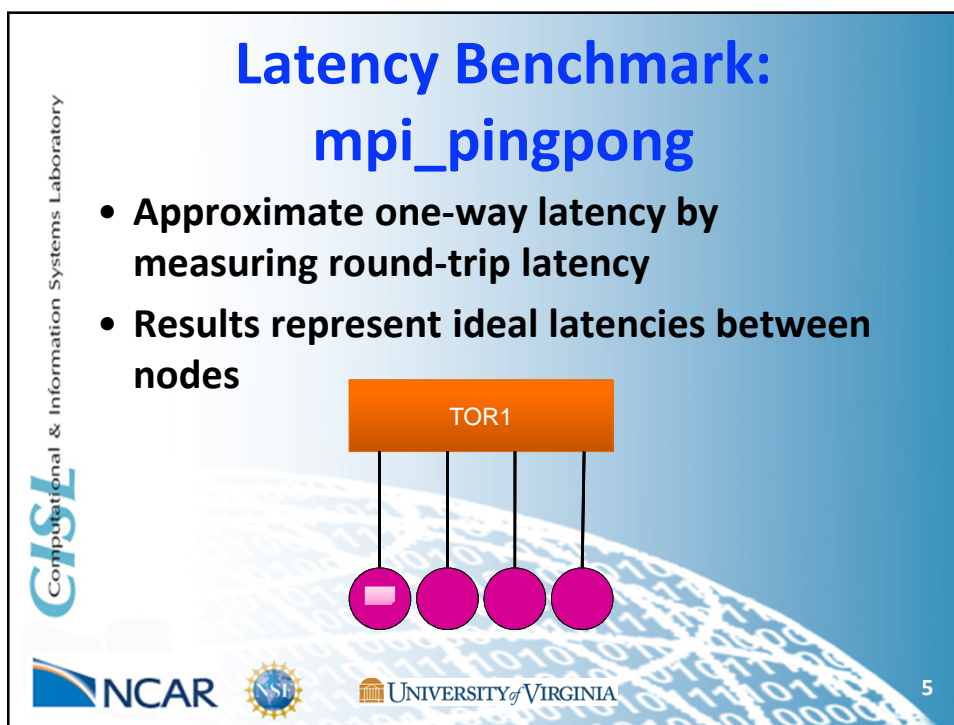
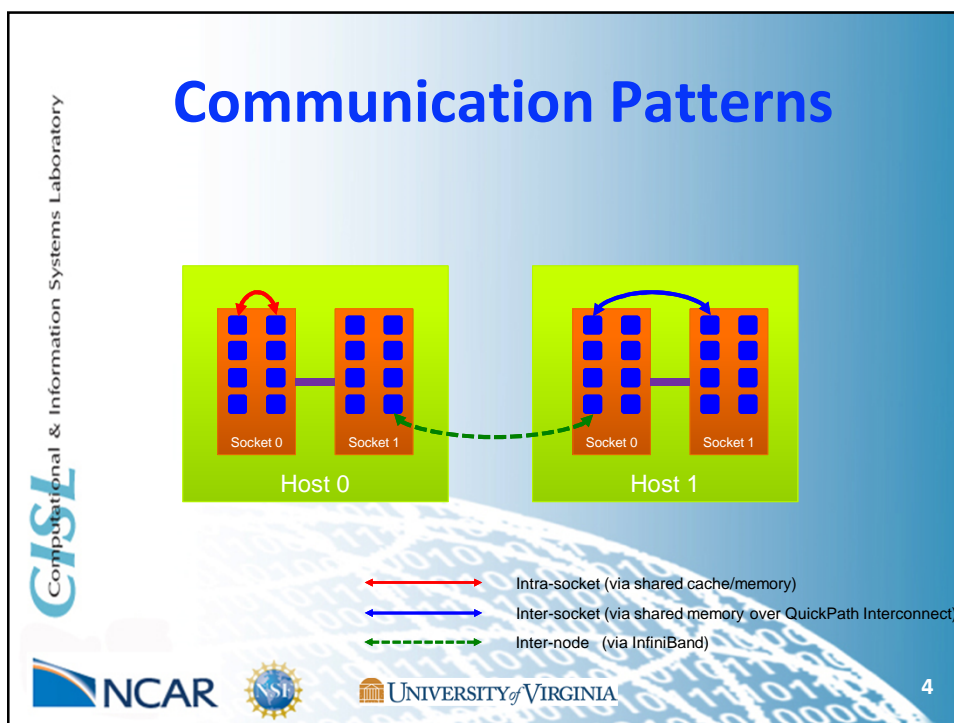
- **Understanding the causes of poor performance of CESM on Yellowstone: a 5-step approach**
 - Experimental execution and data collection
 - CESM trace analysis
 - IBMgtSim: routing study
 - Network simulation
 - Integrated simulation

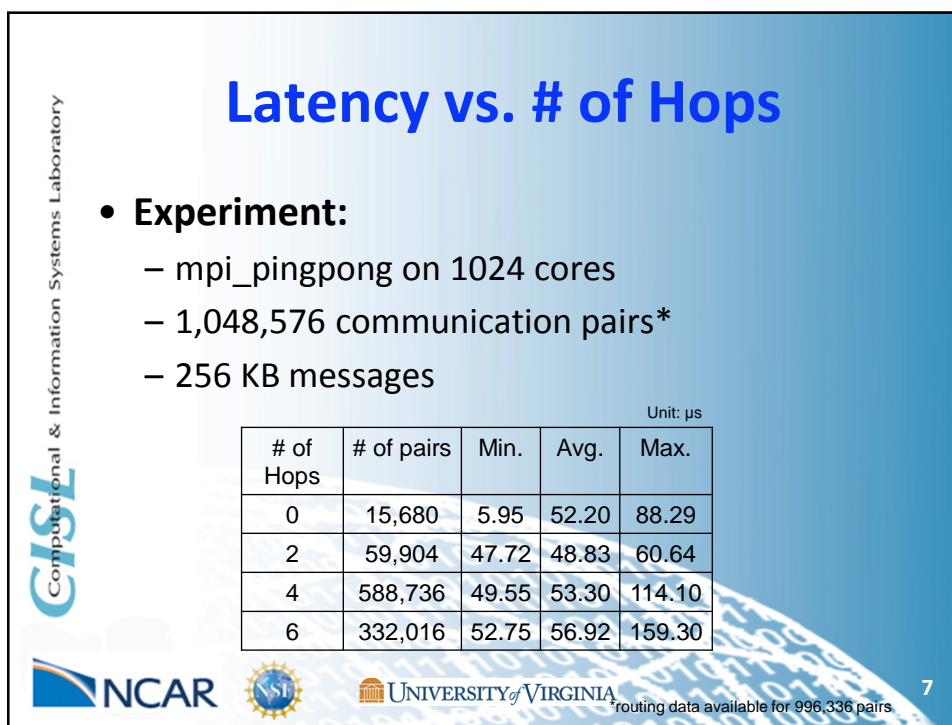
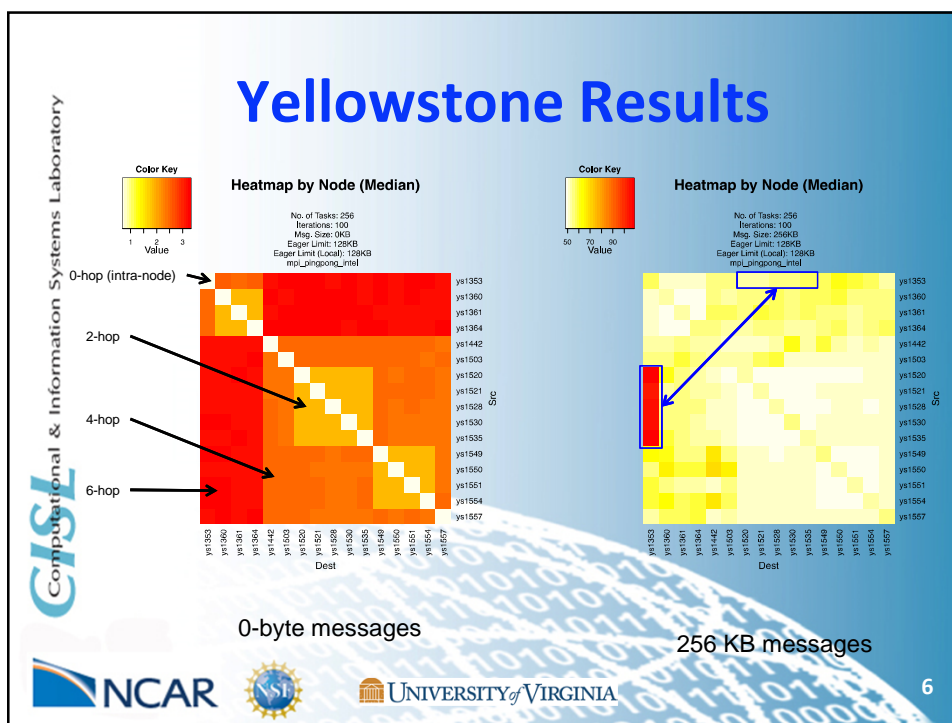
Yellowstone network

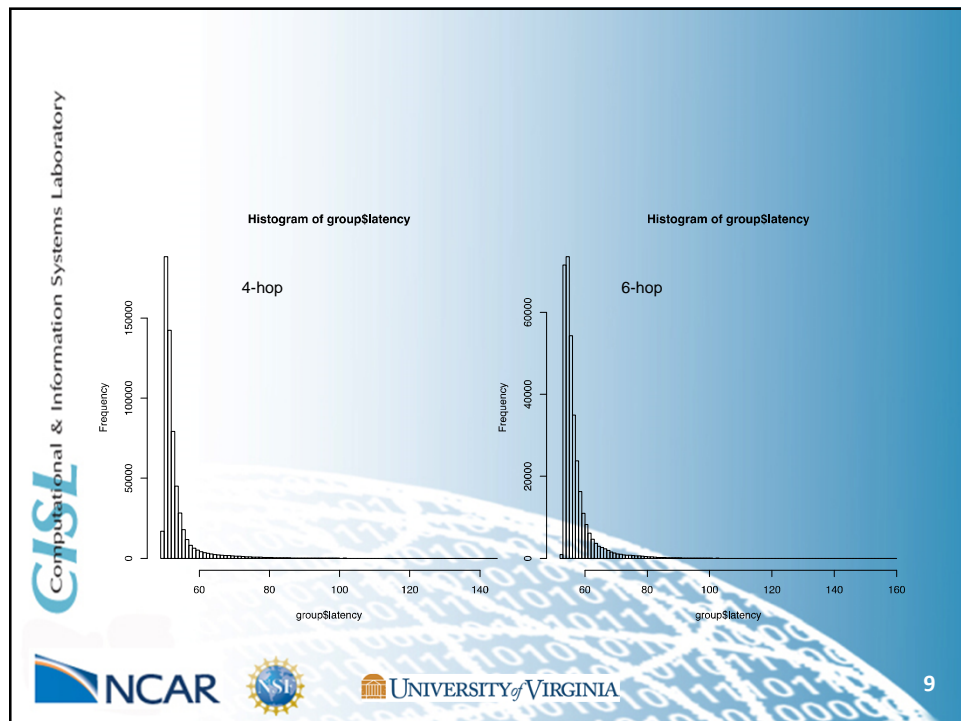
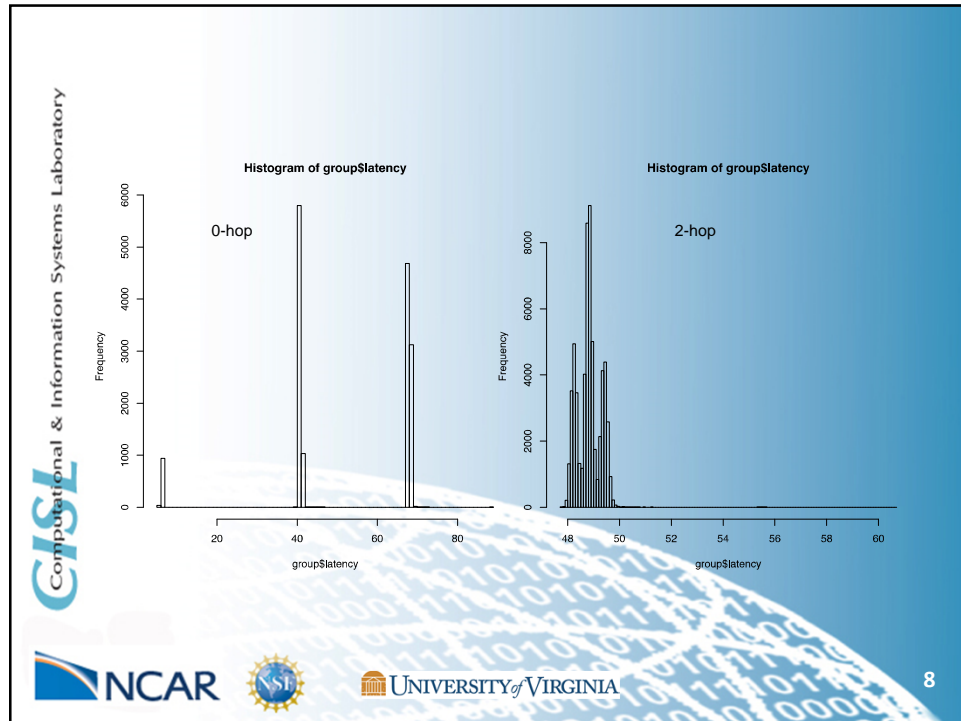


251 A-groups
14 B-groups
9 ORCAs

— 2-hop
— 4-hop
— 6-hop










Computational & Information Systems Laboratory

Bandwidth Benchmark: mpi_bw

- **Measures throughput between two MPI ranks**
- **3 communication patterns:**
 - Intra-socket
 - Inter-socket
 - Inter-node
- **2 communication protocols:**
 - Eager protocol
 - Rendezvous protocol








10

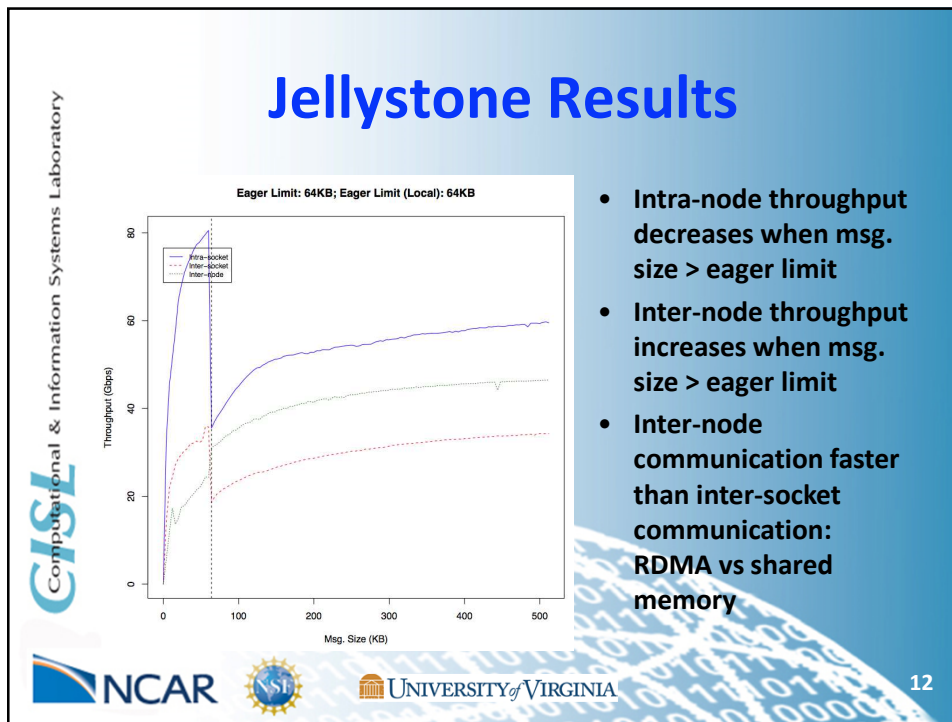
Computational & Information Systems Laboratory

Communication Protocols

- **Rendezvous Protocol: buffer negotiation before sending**
- **Eager Protocol: send directly without confirming available buffer space**
- **InfiniBand: Eager protocol uses SEND/RECV verbs (two-sided communication); Rendezvous protocol uses WRITE/READ verbs (one-sided communication)**
- **Eager Limit: threshold below which Eager protocol is used**

11



Computational & Information Systems Laboratory

IBMgtSim and OpenSM

- **Problem**
 - Special version of OpenSM needed to boot a simulated network
 - Compile time configuration, not supported by Mellanox's version
 - Open source version of OpenSM does not work with Mellanox IBUtils (which contains IBMgtSim)
 - Open source IBUtils outdated
- **Solution**
 - Fix bugs in open source version of OpenSM
 - Patches OpenSM to support FDR link rate
 - Patches to support better simulation work flow (RunSimTest)

NCAR UNIVERSITY of VIRGINIA

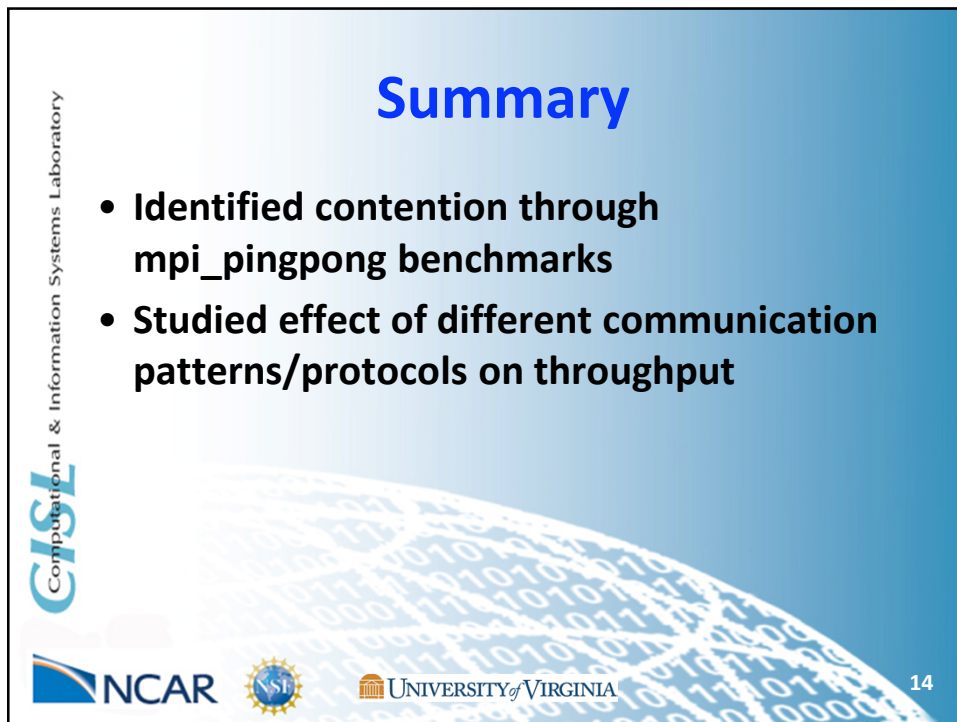
13

Computational & Information Systems Laboratory

Summary

- Identified contention through mpi_pingpong benchmarks
- Studied effect of different communication patterns/protocols on throughput

14



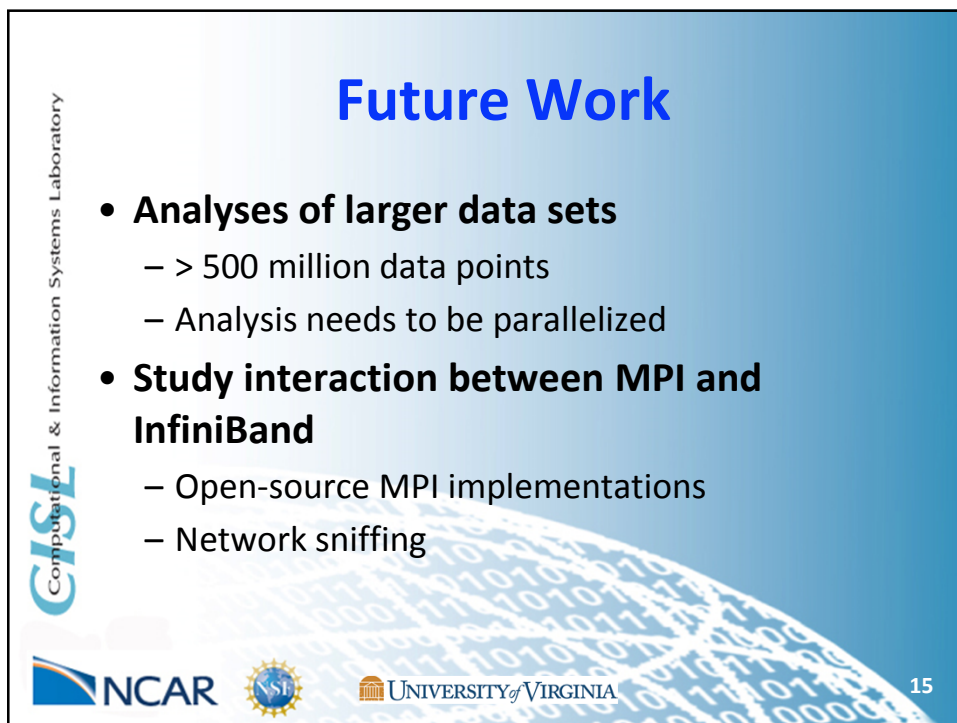
Logos at the bottom: NCAR, NSF, and UNIVERSITY of VIRGINIA.

Computational & Information Systems Laboratory

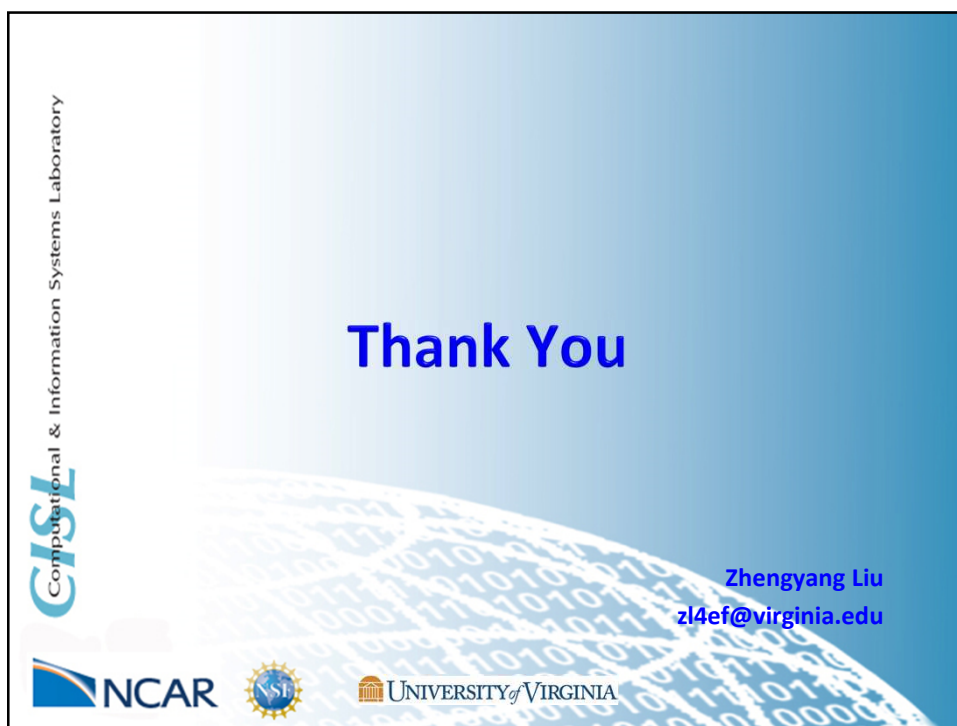
Future Work

- Analyses of larger data sets
 - > 500 million data points
 - Analysis needs to be parallelized
- Study interaction between MPI and InfiniBand
 - Open-source MPI implementations
 - Network sniffing

15



Logos at the bottom: NCAR, NSF, and UNIVERSITY of VIRGINIA.



The slide features a blue gradient background with a white globe graphic at the bottom. The globe is covered in binary code (0s and 1s). The text "Thank You" is centered in a large, bold, blue font. In the bottom left corner, there are three logos: the CIST logo (Computational & Information Systems Laboratory), the NCAR logo, and the University of Virginia logo. In the bottom right corner, the name "Zhengyang Liu" and the email address "zl4ef@virginia.edu" are displayed in blue text.

Computational & Information Systems Laboratory

Thank You

Zhengyang Liu
zl4ef@virginia.edu

NCAR

UNIVERSITY of VIRGINIA